# Introduction to XML and Digital Scholarly Editing Using the Text Encoding Initiative (TEI)

Paul Broyles, North Carolina State University

pabroyle@ncsu.edu

# Workshop Outline

# Overview of Digital Editions

# What is an edition?

- Edition: a particular form in which a work is published; the set of copies of a book produced from a single setting of type.
- Scholarly edition: remediates or represents documents or texts in a form determined by scholarly inquiry.
  - Critical edition: reconstructs a hypothetical text (e.g. a lost copy, an author's original, or an ideal text that never existed) using evidence such as multiple copies, linguistics, etc.
  - Documentary edition: "an edition of a text based on a single document, which attempts to reproduce a certain degree of the peculiarities of the document itself" (Pierazzo, "Digital Documentary Editions and the Others")
- Transcription: a reproduction of the characters that make up a text. *Transcription is interpretation; not neutral.*

# Why digital editions?

- Online publication makes editions easily accessible; wide audience.
- Low publication costs make it economical to edit and publish material that might not be publishable in print; explosion of digital documentary editions.
- High degrees of flexibility and interactivity. No longer necessary to choose between documentary and critical editing; edition can contain both. Innovative forms of representation.
- Ability to tie edited texts to facsimile images of original source material.
- Collaboration and ongoing improvements.
- Makes text available as (structured) data for analysis, visualization.
- Multiple output formats: create encoded digital text once, use in a variety of systems for different purposes.

# Represent primary sources



*Digital Vercelli Book* <http://vbd.humnet.unipi.it/beta2/>, a parallel facsimile and diplomatic edition of an Anglo-Saxon manuscript

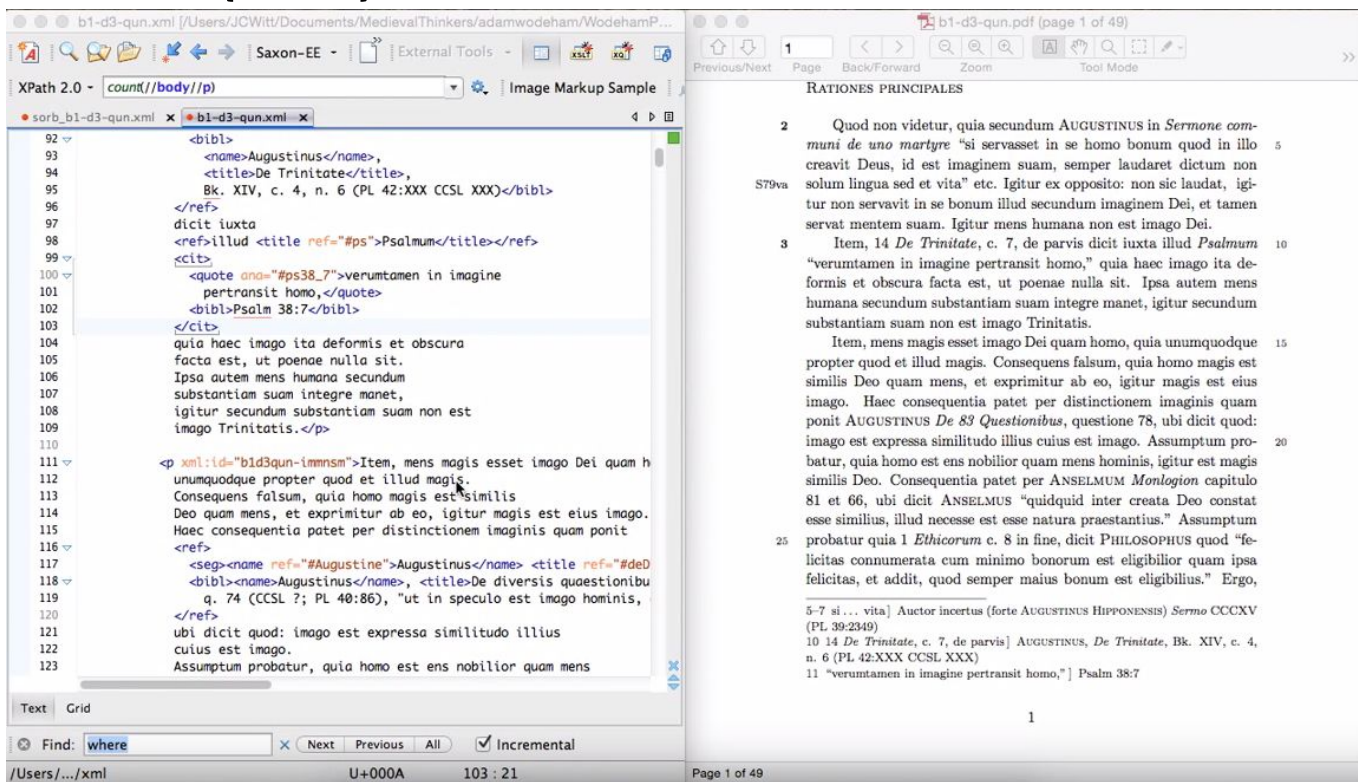# Flexible critical texts



Comparing two manuscript versions of Thoreau's *Walden* with a standard print edition (left). *Walden: A Fluid-Text Edition*. Digital Thoreau. <https://digitalthoreau.org/fluid-text-toc/>.

Digital apparatus for *Piers Plowman: The B-Version Archetype*. Piers Plowman Electronic Archive. <http://piers.chass.ncsu.edu/texts/Bx>.

# Multiple output formats



Demo video for Lbp-Print tool created by Jeffrey Witt for LombardPress <http://lombardpress.org/>.

# Sophisticated analysis and display (maps, graphs, etc.)



From Elisa E. Beshero-Bondar, "Juxtapositions of Place in Thalaba the Destroyer"
<http://ebeshero.github.io/thalaba/index.html>.

# XML and Markup

# What is markup?

- TEI: "we define *markup*, or (synonymously) *encoding*, as any means of making explicit an interpretation of a text" (*Guidelines* v.1).
- Flavors of markup: proofreaders' marks; paragraphing and punctuation marks; critical apparatus; your underlining and notes in the margin.
- McGann: "there is no such thing as an unmarked text, and the markup systems laid upon documents to facilitate computerized analyses are marking orders laid upon already marked up material" (*Radiant Textuality* 138).
- Computer markup: formal codes for embedding information in a document (for example about structure or formatting); facilitates interchange, processing of unfamiliar documents.

# What is XML?

- XML: eXtensible Markup Language
- Privileges the meaning/structure of information over presentation (how it looks)
- Metalanguage: describes other languages that can be defined for specific purposes (hence *extensible*: not limited to preset tags)
- Separates content from metadata (at least in principle)
- Provides a mechanism for *validation*: making sure the document follows the rules
- Allows for the exchange of a wide variety of data between different systems

# Tags, Elements, and Attributes



An XML element with no content is called an **empty element** and can be expressed with a single **self-closing tag**: `<tag attribute="value"/>`

# XML as Ordered Hierarchy

book
- introduction
- chapter
  - heading
  - section
  - section
- chapter
  - heading
  - section
  - section
- index

```
1   <?xml version="1.0" encoding="UTF-8"?>
2   <book>
3       <introduction/>
4       <chapter>
5           <heading/>
6           <section/>
7           <section/>
8       </chapter>
9       <chapter>
10          <heading/>
11          <section/>
12          <section/>
13      </chapter>
14      <index/>
15  </book>
```

Elements **nest** inside each other—that is, one element can contain another element. All the elements in an XML document are contained within the **root element**. The overall structure of an XML document is that of a **tree**.

Source: Michelle Dalmau and John Walsh, "Text Encoding Initiative Workshop: Intro to Text Encoding," Indiana University <http://dcl.slis.indiana.edu/teiworkshop/teiworkshop_i.pdf>.

# The TEI

# What is the TEI?

- TEI = Text Encoding Initiative
- TEI is an organization, the TEI Consortium, which develops standards for digital representation of texts. Founded 1987.
- TEI Guidelines define an encoding language, currently expressed in terms of XML. (When I refer to the TEI, I typically mean this encoding format.)

# What does TEI encoding do?

- Allows (primarily textual) documents to be encoded with a high degree of descriptive precision.
  - Not just a transcription: appearance of the document, logical structure, information contained within it, etc.
- Provides a set of elements and attributes that can be used to describe the components of documents, as well as rules for *validating* documents.
- Shared standards enable exchange, interoperability.
- Suite of tools that make it simple to transform and publish TEI documents in multiple formats.
- Flexible and customizable: can define specific, documented variations to meet the need of particular projects.

# General Approaches to TEI Encoding

- Organize documents around their logical structure: division into meaningful sections, hierarchies.
    - For example, hierarchy of a book likely to be organized around chapters and paragraphs rather than pages and lines.
    - **But** there are also approaches that prioritize arrangement on the page; flexible to needs of project, material.
- Encode/describe document as it exists, not as you want it to appear.
    - Precise mechanisms for describing formatting of input; formatting of output is left to mechanisms beyond the TEI file.
- Use as much or as little of the TEI tag set as fits the project.
- There are many equally valid ways to encode the same document depending on your interpretation, what features you want to emphasize/analyze.

# Hands-On With TEI

# TEI Web Editor

- Web app developed by Jeffrey Witt (Assistant Professor of Philosophy, Loyola University Maryland)
- Integrated with GitHub (online saving), Mirador IIIF viewer (show images in editing window)
- Checks if your work is well-formed (does it follow the syntax of XML), but **not** if it validates (does it follow the rules of TEI)
- Live preview of your work, lightly formatted
- Online demo (may not be permanent) at https://tei-web-editor.herokuapp.com.

# Workshop Instructions

For instructions on completing the workshop exercises, visit
https://paulbroyles.github.io/tei-workshop-instructions/.

# Where do we go from here?

# Learning about TEI

- TEI by Example (http://teibyexample.org/)
- TEI-L Listserv – best for specific questions after you've consulted other resources (https://listserv.brown.edu/archives/cgi-bin/wa?A0=TEI-L)
- The TEI Guidelines (http://www.tei-c.org/release/doc/tei-p5-doc/en/html/)
- Consult with me – pabroyle@ncsu.edu

# Reading the TEI Guidelines

- Organized into chapters describing modules that serve a specific purpose (usually, a type of material or approach to encoding).
- Chapters are long and detailed; use the table of contents (and search) to find what you want. You don't need to read everything!
- Appendix C, Elements, and D, Attributes, are especially useful for getting definitions of particular elements and figuring out where they can go and what they can do.

# Software for XML and TEI

- XML is a text-based language, can be edited in any text editor (e.g. Atom, BBEdit [Mac], Notepad [Windows]).
- Dedicated software for editing XML can help avoid errors by checking that XML is well-formed (correct syntax) and valid (follows the rules defined by TEI).
- Commonly used in TEI community: oXygen (http://oxygenxml.com/). Commercial software. Available in CHASS labs.
- oXygen includes built-in support for TEI: provides template for new TEI documents. Also built-in support for transformations with XSLT.

# Publishing TEI Editions

- The simple way: CETEIcean
  - What TEI Web Editor uses to publish your work on GitHub Pages
  - JavaScript library to make TEI documents displayable in web browser
  - Create very simple HTML page referencing both the TEI document and CETEIcean JS
  - Can build more sophisticated experiences with CSS to style, additional JavaScript
- The old way: TEI Boilerplate
  - Similar to CETEIcean, but uses a language called XSLT to transform TEI into HTML in the web broswer
  - May not have a long-term future and does not support current versions of XSLT
  - Today, CETEIcean almost always a better choice

# Publishing TEI Editions: More Sophisticated Tools

- TEI Publisher
  - Tool for building sophisticated web apps for TEI editions
  - Easy support for multiple documents, search, navigation, etc.
  - Deep formatting control, multiple outputs (PDF, ePub, etc.)
  - Much more complicated setup than CETEIcean, especially for single document; requires more advanced hosting (runs on eXist database).
- EVT (Edition Visualization Technology)
  - Tool for displaying diplomatic editions in parallel with facsimile images
  - Simpler to set up, host than TEI Publisher, but still much more complicated than CETEIcean
  - Main support for diplomatic editions, though a beta version promises support for critical editions.

# Analyzing TEI Texts

- Juxta Commons (text collation tool) accepts TEI XML as both an input and an output format; can aid in the editing process or be used to analyze differences among witnesses/versions of a text
- Voyant Tools (text analysis tool) is TEI-aware, will allow you to upload a file in TEI or other XML vocabulary and identify title, text, etc.
- Because TEI files consist of structured data, can be processed to use with any software. Marking up features of interest and providing metadata means you can extract things you can't from plaintext: frequency of correction, how often a person is referenced even if not by name, etc.

# Publishing and Transforming TEI: XSLT

- XSLT (eXtensible Stylesheet Language Transformations): a language for transforming XML documents (itself written as a form of XML)
- Can be used to make systematic changes to XML, or to convert to another format (HTML, JSON)
- Many TEI projects use XSLT on the server to create custom interfaces
- XSLT (and the related language XQuery) can be used to process TEI files, gathering data from them directly or turning into another format to use with text processing software
- Sophisticated and complicated programming language, not needed for basic publishing, but worth learning if you want to do deep work